

707.000
Web Science and Web Technology
„Network Theory and Terminology“

Markus Strohmaier

Univ. Ass. / Assistant Professor
Knowledge Management Institute
Graz University of Technology, Austria

e-mail: markus.strohmaier@tugraz.at
web: <http://www.kmi.tugraz.at/staff/markus>

Overview

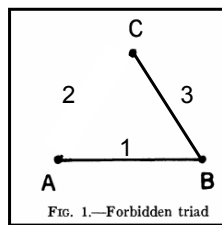
Agenda

- A selection of *relevant* concepts from Graph and Network Theory

Bridges and Strong Ties [Granovetter 1973]

Example:

1. Imagine the strong tie between A and B
2. Imagine the strong tie between B and C
3. Then, the forbidden triad **implies** that a tie **exists** between C and A
(it forbids that a tie between C and A does not exist)
1. From that follows, that A-B is not a bridge (because there is another path A-B that goes through C)



Why is this interesting?

⇒ Strong ties can be a bridge ONLY IF neither party to it has any other strong ties

⇒ Highly unlikely in a social network of any size

⇒ Weak ties suffer no such restriction, though they are not automatically bridges

⇒ But, **all bridges are weak ties**

In Reality [Granovetter 1973]

it probably happens only rarely, that a specific tie provides the only path between two points

Local bridges: the shortest path between its two points (other than itself)

- Bridges are efficient paths
- Alternatives are more costly
- Local bridges of degree n
- A local bridge is more significant as its degree increases

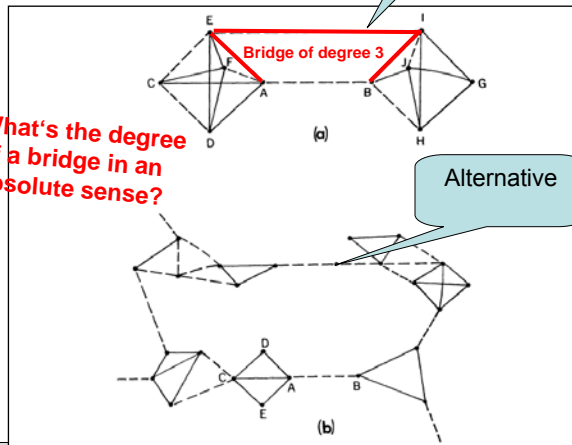
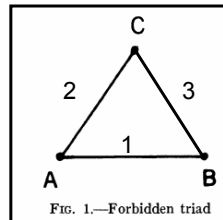


FIG. 2.—Local bridges. a, Degree 3; b, Degree 13. — = strong tie; - - - = weak tie.

In Reality ...

Strong ties can represent *local* bridges BUT
They are weak (i.e. they have a low degree)

Why?



What's the degree of the local bridge A-B?

Implications of Weak Ties [Granovetter 1973]

- Those weak ties, that are local bridges, create more, and shorter paths.
- The removal of the average weak tie would do more damage to transmission probabilities than would that of the average strong one
- **Paradox:** While *weak ties* have been denounced as generative of alienation, *strong ties*, breeding local cohesion, lead to overall fragmentation

What are sources of weak ties/bridges?

Can you identify some implications for social networks on the web / for search in these networks?

How does this relate to Milgram's experiment?

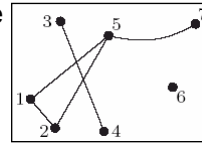
Completion rates in Milgram's experiment were reported higher for acquaintance than friend relationships [Granovetter 1973]

Terminology

<http://www.cis.upenn.edu/~Emkearns/teaching/NetworkedLife/>
 [Diestel 2005]

Network

- A collection of individual or atomic entities
- Referred to as nodes or vertices (the “dots” or “points”)
- Collection of links or edges between vertices (the “lines”)
- Links can represent any pairwise relationship
- Links can be directed or undirected
- Network: entire collection of nodes and links
- For us, a network is an abstract object (list of pairs) and is separate from its visual layout
- that is, we will be interested in properties that are invariant
 - structural properties
 - statistical properties of families of networks



Social Networks

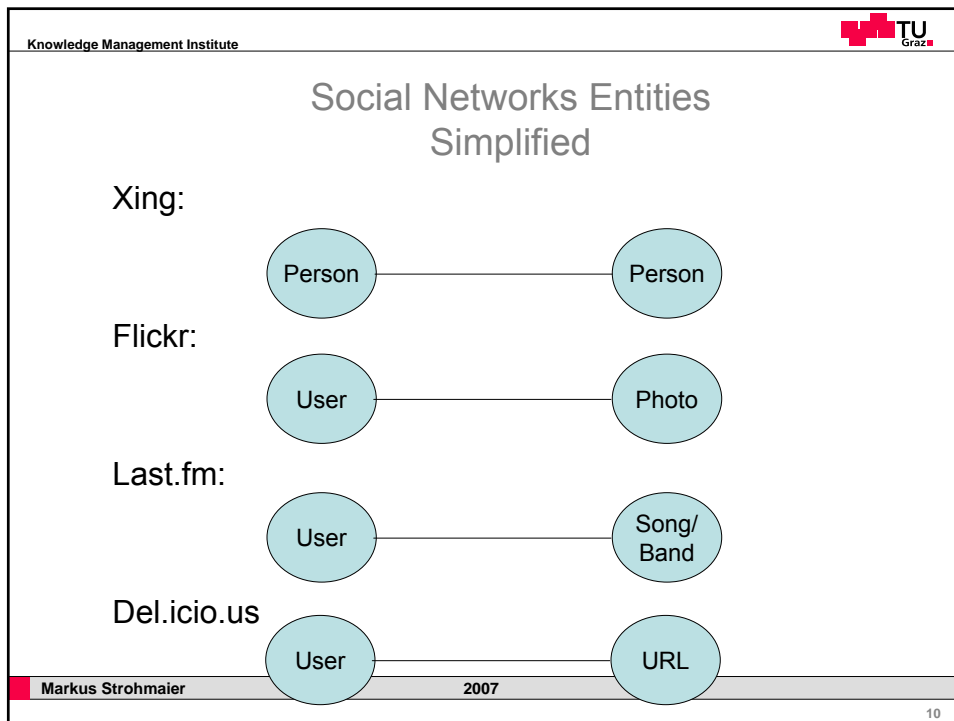


Figure 1.3. Real social networks exhibit clustering, the tendency of two individuals who share a mutual friend to be friends themselves. Here, Ego has six friends, each of whom is friends with at least one other.

Knowledge Management Institute TU Graz

Social Networks Examples

The screenshot shows the del.icio.us interface. At the top, there are navigation links for 'popular' and 'recent', along with 'login', 'register', and 'help'. A search bar contains the URL 'http://www.devhardware.com/c/e/'. Below the search bar, the main content area displays the title 'Why and How to Flash Your BIOS' and a note that 'this url has been saved by 106 people'. A 'user notes' section shows a note from 'r1aw77' dated 'Aug '07'. On the right side, there are sections for 'common tags' (including 'bios', 'hardware', 'tutorial') and 'posting history'.



Network Examples [Newman 2003]

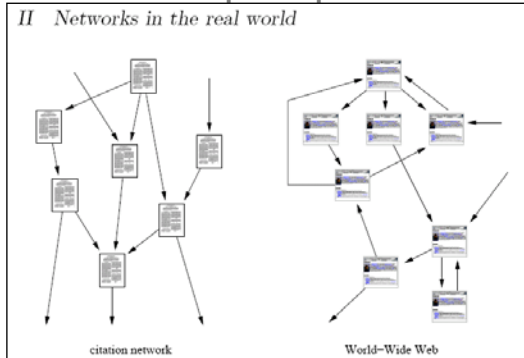


FIG. 4 The two best studied information networks. Left: the citation network of academic papers in which the vertices are papers and the directed edges are citations of one paper by another. Since papers can only cite those that came before them (lower down in the figure) the graph is acyclic—it has no closed loops. Right: the World Wide Web, a network of text pages accessible over the Internet, in which the vertices are pages and the directed edges are hyperlinks. There are no constraints on the Web that forbid cycles and hence it is in general cyclic.

Terminology II

<http://www.cis.upenn.edu/~Emkearns/teaching/NetworkedLife/>

- Network size: total number of vertices (denoted N)
- Maximum number of edges (undirected): $N(N-1)/2 \sim N^2/2$
- Distance or geodesic path between vertices u and v :
 - number of edges on the **shortest path** from u to v
 - can consider directed or undirected cases
 - infinite if there is no path from u to v
- Diameter of a network
 - worst-case diameter: largest distance between a pair
 - Diameter: longest shortest path between any two pairs
 - average-case diameter: average distance
- If the distance between all pairs is finite, we say the network is connected; else it has multiple components
- Degree of vertex v : number of edges connected to v
- Density: ratio of edges to vertices

Definitions

[Newman 2003]

Vertex (pl. vertices): The fundamental unit of a network, also called a site (physics), a node (computer science), or an actor (sociology).

Edge: The line connecting two vertices. Also called a bond (physics), a link (computer science), or a tie (sociology).

Directed/undirected: An edge is directed if it runs in only one direction (such as a one-way road between two points), and undirected if it runs in both directions. Directed edges, which are sometimes called *arcs*, can be thought of as sporting arrows indicating their orientation. A graph is directed if all of its edges are directed. An undirected graph can be represented by a directed one having two edges between each pair of connected vertices, one in each direction.

Degree: The number of edges connected to a vertex. Note that the degree is not necessarily equal to the number of vertices adjacent to a vertex, since there may be more than one edge between any two vertices. In a few recent articles, the degree is referred to as the "connectivity" of a vertex, but we avoid this usage because the word connectivity already has another meaning in graph theory. A directed graph has both an in-degree and an out-degree for each vertex, which are the numbers of in-coming and out-going edges respectively.

Component: The component to which a vertex belongs is that set of vertices that can be reached from it by paths running along edges of the graph. In a directed graph a vertex has both an in-component and an out-component, which are the sets of vertices from which the vertex can be reached and which can be reached from it.

Geodesic path: A geodesic path is the shortest path through the network from one vertex to another. Note that there may be and often is more than one geodesic path between two vertices.

Diameter: The diameter of a network is the length (in number of edges) of the longest geodesic path between any two vertices. A few authors have also used this term to mean the average geodesic distance in a graph, although strictly the two quantities are quite distinct.

Terminology III

<http://www.infosci.cornell.edu/courses/info204/2007sp/>

[Diestel 2005]

In undirected networks

- Paths

- A sequence of nodes $v_1, \dots, v_i, v_{i+1}, \dots, v_k$ with the property that each consecutive pair v_i, v_{i+1} is joined by an edge in G

- Cycles (in undirected networks)

- A path with $v_1 = v_k$ (Begin and end node are the same)
 - Cyclic vs. Acyclic (not containing any cycles: e.g. forests) networks

In directed networks

- Path or cycles must respect directionality of

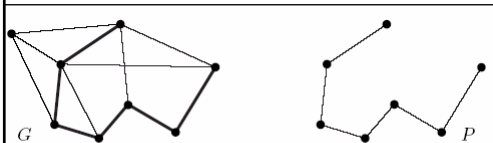


Fig. 1.3.1. A path $P = P^6$ in G

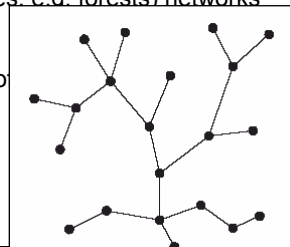


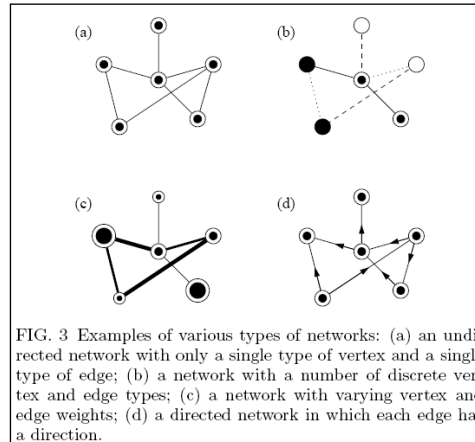
Fig. 1.5.1. A tree

Examples

[Newman 2003]

Undirected,
single edge and
node type

Undirected,
varying edge and
node weights



Undirected,
multiple edge
and node types

Directed, each
edge has a
direction

Terminology IV

<http://www.infosci.cornell.edu/courses/info204/2007sp/>

- **Average Pairwise Distance**
 - The average distance between all pairs of nodes in a graph. If the graph is unconnected, the average distance between all pairs in the largest component.
- **Connectivity**
 - An undirected graph is connected if for every pair of nodes u and v , there is a path from u to v (there is not more than one component).
 - A directed graph is strongly connected if for every two nodes u and v , there is a path from u to v and a path from v to u
- **Giant Component**
 - A single connected component that accounts for a significant fraction of all nodes

Average degree k

<http://www.infosci.cornell.edu/courses/info204/2007sp/>

- Average degree k
 - Degree: The number of edges for which a node is an endpoint
 - In undirected graphs: number of edges
 - In directed graphs: k_{in} and k_{out}
 - Average degree: average of the degree of all nodes, a measure for the density of a graph

$$d(G) := \frac{1}{|V|} \sum_{v \in V} d(v)$$

Degree Distributions

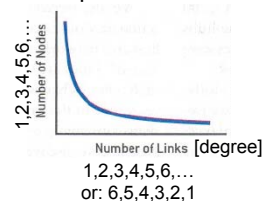
[Barabasi and Bonabeau 2003]

- Degree distribution $p(k)$
 - A plot showing the fraction of nodes in the graph of degree k , for each value of k

Related concepts

- Degree histogram
- Rank / frequency plot
- Cumulative Degree function (CDF)
- Pareto distribution

Example:



Degree Distributions Examples

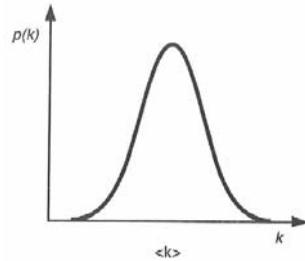


Figure 4.1. The normal distribution specifies the probability, $p(k)$, that a randomly selected node will have k neighbors. The average degree $\langle k \rangle$ lies at the peak of the distribution.

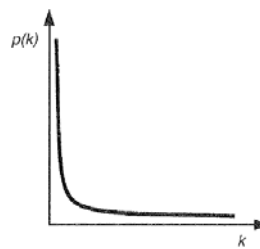


Figure 4.2. A power-law distribution. Although it decreases rapidly with k , it does so much slower than the normal distribution in figure 4.1, implying that large values of k are more likely.

Clustering Coefficient

<http://www.infosci.cornell.edu/courses/info204/2007sp/>

- Clustering Coefficient C
 - Triangles or closed triads: Three nodes with edges between all of them
 - over all sets of three nodes in the graph that form a connected set (i.e. one of the three nodes is connected to all the others), what fraction of these sets in fact form a triangle?
 - This fraction can range from 0 (when there are no triangles) to 1 (for example, in a graph where there is an edge between each pair of nodes — such a graph is called a clique, or a complete graph).
 - Or in other words: The clustering coefficient gives the fraction of pairs of neighbors of a vertex that are adjacent, averaged over all vertices of the graph. [p344, Brandes and Erlebach 2005]
 - Page 88, [Watts 2005]
 - Related: „Transitivity“

Clustering Coefficient

Images taken from http://en.wikipedia.org/w/index.php?title=Clustering_coefficient&oldid=152650779

- Number of edges between neighbours of a **given node** divided by the number of possible edges between neighbours

- Directed Graphs

$$C_i = \frac{|\{e_{jk}\}|}{k_i(k_i - 1)} : v_j, v_k \in N_i, e_{jk} \in E.$$

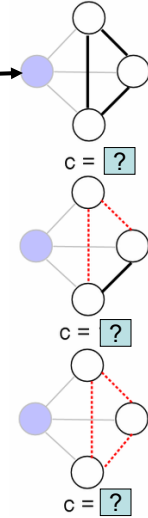
- Undirected Graphs

$$C_i = \frac{2|\{e_{jk}\}|}{k_i(k_i - 1)} : v_j, v_k \in N_i, e_{ij} \in E.$$

Degree

Edges between neighbourhood nodes

Neighbourhood nodes



Graph Theory & Network Theory

- Graph Theory

- Mathematics of graphs
- Networks with pure structure with properties that are fixed over time
- Focus on syntax rather than semantics
 - Nodes and edges do not have semantics
 - E.g. A node does not have a social identity
- Concerned with characteristics of graphs
- Proofs
- Algorithms

- Network Theory

- Relate to real-world phenomena
 - Social networks
 - Economic networks
 - Energy networks
- Networks are *doing something*
 - *Making new relations*
 - *Making money*
 - *Producing power*
- Are dynamic
 - Structure: Dynamics of the network
 - Agency: Dynamics in the network
- Are active, which effects
 - *Individual behavior*
 - *Behavior of the network as a whole*

Networks [Watts 2003]

TABLE 3.2 STATISTICS OF SMALL WORLD NETWORKS

	L_{ACTUAL}	L_{RANDOM}	C_{ACTUAL}	C_{RANDOM}
MOVIE ACTORS	3.65	2.99	0.79	0.00027
POWER GRID	18.7	12.4	0.080	0.005
<i>C. ELEGANS</i>	2.65	2.25	0.28	0.05

L =Path Length; C =Clustering Coefficient.

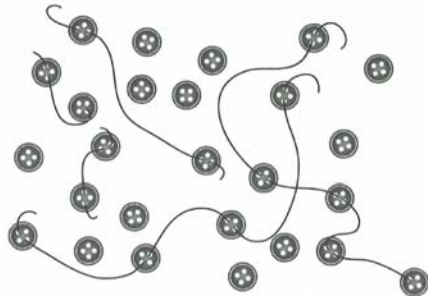
Compared to
imaginary random
networks

Network Theory

- Are there general statements we can make about *any* class of network?
- A Science of Networks

Random Networks

- Page 44/ff, Watts 2003, random graphs



Random graph: a network of nodes connected by links in a purely random fashion.

Analogy of Stuart Kaufmann: Throw a boxload of buttons onto the floor, then choose pairs of buttons at random tying them together

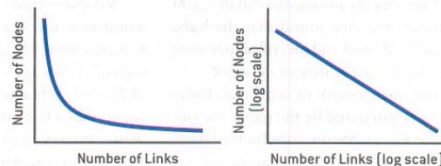
Figure 2.1. A random graph imagined as a collection of buttons tied by strings. Pairs of nodes (buttons) are connected at random by links or ties.

Scale-Free Networks

[Barabasi and Bonabeau 2003]

- Some nodes have a tremendous number of connections to other nodes (hubs), whereas most nodes have just a handful
- Robust against accidental failures, but vulnerable to coordinated attacks
- Popular nodes can have millions of links: The network appears to have no scale (no limit)
- Two prerequisites: [watts2003]
 - Growth
 - Preferential attachment
- Problem:
 - Scale-free networks are only ever truly scale-free when the network is infinitely large (whereas in practice, they are mostly not)
 - This introduces a cut off [page 111, watts 2003]

Power Law Distribution of Node Linkages



Scale-free Networks

[Watts 2003]

The alpha parameter

- $y = C x^{-\alpha}$ (C, α being constants) or $\log(y) = \log(C) - \alpha \log(x)$
- a power-law with exponent α is depicted as a straight line with slope $-\alpha$ on a log-log plot

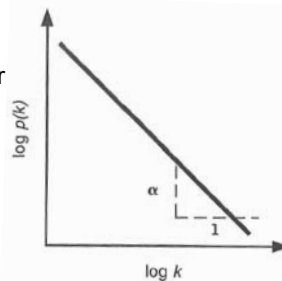


Figure 4.3. A power-law distribution on a log-log plot. The exponent α is the slope of the line (the line drops by α for each unit on the horizontal axis).

Examples

- If a number of cities of a given size decreases in inverse proportion to the size, then we say the distribution has an exponent of [one/two]

That means, we are likely to see cities such as Graz (250.000) roughly [ten/hundred] times as frequently as cities like Vienna (including the Greater Vienna Area that is roughly 10 times larger)

Networks [Newman 2003]

	network	type	n	m	z	ℓ	α	$C^{(1)}$	$C^{(2)}$	r	Ref(s).
social	film actors	undirected	449913	25 516 482	113.43	3.48	2.3	0.20	0.78	0.208	20, 416
	company directors	undirected	7 673	55 392	14.44	4.60	-	0.59	0.88	0.276	105, 323
	math coauthorship	undirected	253 339	496 489	3.92	7.57	-	0.15	0.34	0.120	107, 182
	physics coauthorship	undirected	52 909	245 300	9.27	6.19	-	0.45	0.56	0.363	311, 313
	biology coauthorship	undirected	1 520 251	11 803 064	15.53	4.92	-	0.088	0.60	0.127	311, 313
	telephone call graph	undirected	47 000 000	80 000 000	3.16	-	2.1	-	-	-	8, 9
	email messages	directed	59 912	86 300	1.44	4.95	1.5/2.0	-	0.16	-	136
	email address books	directed	16 881	57 029	3.38	5.22	-	0.17	0.13	0.092	321
	student relationships	undirected	573	477	1.66	16.01	-	0.005	0.001	-0.029	45
sexual contacts	undirected	2 810	-	-	-	3.2	-	-	-	265, 266	
information	WWW nd.edu	directed	269 504	1 497 135	5.55	11.27	2.1/2.4	0.11	0.29	-0.067	14, 34
	WWW Altavista	directed	203 549 046	2 130 000 000	10.46	16.18	2.1/2.7	-	-	-	74
	citation network	directed	783 339	6 716 198	8.57	-	3.0/-	-	-	-	351
	Roget's Thesaurus	directed	1 022	5 103	4.99	4.87	-	0.13	0.15	0.157	244
word co-occurrence	undirected	460 902	17 000 000	70.13	-	2.7	-	0.44	-	119, 157	
technological	Internet	undirected	10 697	31 992	5.98	3.31	2.5	0.035	0.39	-0.189	86, 148
	power grid	undirected	4 941	6 594	2.67	18.99	-	0.10	0.080	-0.003	416
	train routes	undirected	587	19 603	66.79	2.16	-	0.69	-	-0.033	366
	software packages	directed	1 439	1 723	1.20	2.42	1.6/1.4	0.070	0.082	-0.016	318
	software classes	directed	1 377	2 213	1.61	1.51	-	0.033	0.012	-0.119	395
	electronic circuits	undirected	24 097	53 248	4.34	11.05	3.0	0.010	0.030	-0.154	155
	peer-to-peer network	undirected	880	1 296	1.47	4.28	2.1	0.012	0.011	-0.306	6, 354
biological	metabolic network	undirected	765	3 686	9.64	2.56	2.2	0.090	0.67	-0.240	214
	protein interactions	undirected	2 115	2 240	2.12	6.80	2.4	0.072	0.071	-0.156	212
	marine food web	directed	135	598	4.43	2.05	-	0.16	0.23	-0.263	204
	freshwater food web	directed	92	997	10.84	1.90	-	0.20	0.087	-0.326	272
	neural network	directed	307	2 359	7.68	3.97	-	0.18	0.28	-0.226	416, 421

TABLE II Basic statistics for a number of published networks. The properties measured are: type of graph, directed or undirected; total number of vertices n ; total number of edges m ; mean degree z ; mean vertex-vertex distance ℓ ; exponent α of degree distribution if the distribution follows a power law (or "-" if not; in/out-degree exponents are given for directed graphs); clustering coefficient $C^{(1)}$ from Eq. (3); clustering coefficient $C^{(2)}$ from Eq. (6); and degree correlation coefficient r , Sec. III.F. The last column gives the citation(s) for the network in the bibliography. Blank entries indicate unavailable data.

Scale-Free Networks

- cut off [page 111, watts 2003]

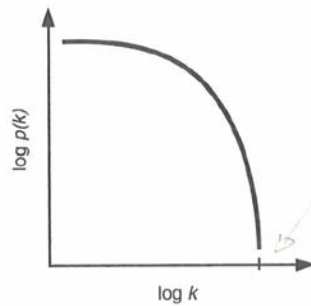


Figure 4.4. A normal-type distribution on a log-log plot. The *cutoff* occurs where the curve disappears into the horizontal axis.

no afa
log. size
in...

Scale-Free Networks

- cut off [page 111, watts 2003]

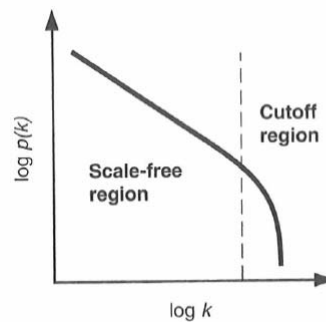


Figure 4.5. In practice, power-law distributions always display a characteristic cutoff because of the finite size of the system. The observed degree distribution, therefore, is only ever a straight line on a log-log plot, over some range.

Limited maximum degree because of the finite set of nodes in a network

Examples of Scale-Free Networks

[Newman 2003]

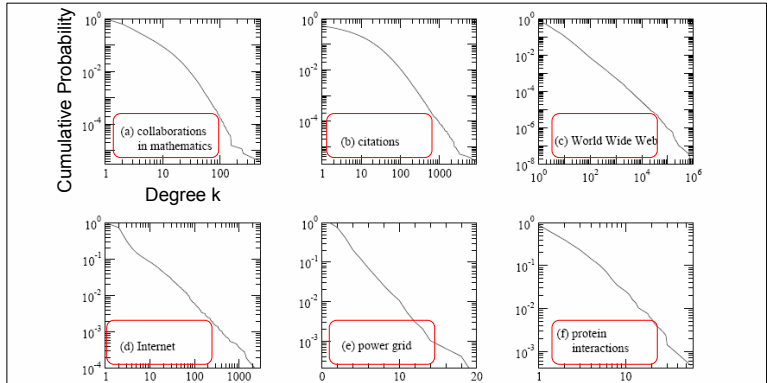


FIG. 6 Cumulative degree distributions for six different networks. The horizontal axis for each panel is vertex degree k (or in-degree for the citation and Web networks, which are directed) and the vertical axis is the cumulative probability distribution of degrees, i.e., the fraction of vertices that have degree greater than or equal to k . The networks shown are: (a) the collaboration network of mathematicians [182]; (b) citations between 1981 and 1997 to all papers cataloged by the Institute for Scientific Information [351]; (c) a 300 million vertex subset of the World Wide Web, circa 1999 [74]; (d) the Internet at the level of autonomous systems, April 1999 [86]; (e) the power grid of the western United States [416]; (f) the interaction network of proteins in the metabolism of the yeast *S. Cerevisiae* [212]. Of these networks, three of them, (c), (d) and (f), appear to have power-law degree distributions, as indicated by their approximately straight-line forms on the doubly logarithmic scales, and one (b) has a power-law tail but deviates markedly from power-law behavior for small degree. Network (e) has an exponential degree distribution (note the log-linear scales used in this panel) and network (a) appears to have a truncated power-law distribution of some type, or possibly two separate power-law regimes with different exponents.

Graph Structure in the Web

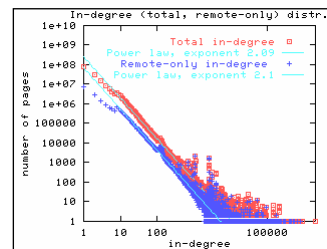
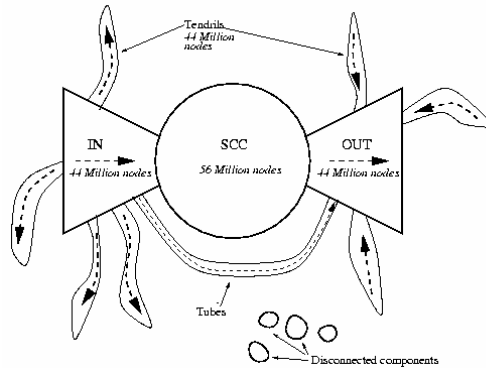
[Broder et al 2000]

Most (over 90%) of the approximately 203 million nodes in a May 1999 crawl form a connected component if links are treated as *undirected* edges.

IN consists of pages that can reach the SCC, but cannot be reached from it

OUT consists of pages that are accessible from the SCC, but do not link back to it

TENDRILS contain pages that cannot reach the SCC, and cannot be reached from the SCC



Interesting Results

[Broder et al 2000]

- the diameter of the central core (SCC) is at least 28, and that the diameter of the graph as a whole is over 500
- for randomly chosen source and destination pages, the probability that any path exists from the source to the destination is only 24%
- if a directed path exists, its average length will be about 16
- if an undirected path exists (i.e., links can be followed forwards or backwards), its average length will be about 6

Scale-Free vs. Random Networks

[Barabasi and Bonabeau 2003]

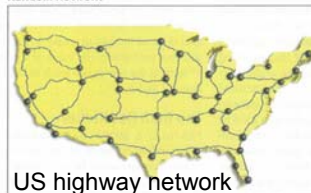
RANDOM VERSUS SCALE-FREE NETWORKS

RANDOM NETWORKS, which resemble the U.S. highway system (simplified in left map), consist of nodes with randomly placed connections. In such systems, a plot of the distribution of node linkages will follow a bell-shaped curve (left graph), with most nodes having approximately the same number of links.

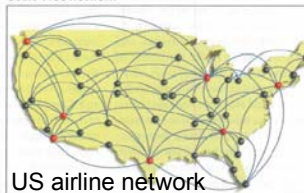
In contrast, scale-free networks, which resemble the U.S. airline system (simplified in right map), contain hubs (red)—

nodes with a very high number of links. In such networks, the distribution of node linkages follows a power law (center graph) in that most nodes have just a few connections and some have a tremendous number of links. In that sense, the system has no "scale." The defining characteristic of such networks is that the distribution of links, if plotted on a double-logarithmic scale (right graph), results in a straight line.

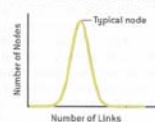
Random Network



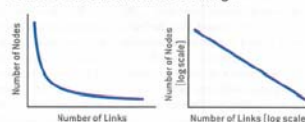
Scale-Free Network



Bell Curve Distribution of Node Linkages



Power Law Distribution of Node Linkages



Bipartite Networks [Watts 2003]

- Page 120
- Can always be represented as unipartite networks

Can you give examples for bipartite networks on the web?

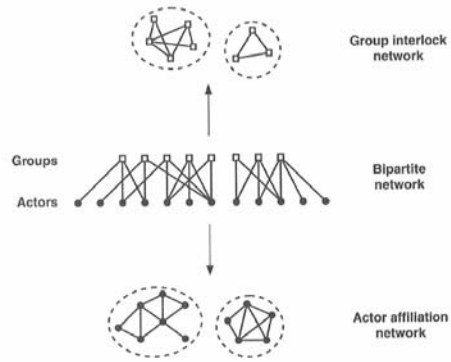


Figure 4.6. Affiliation networks are best represented as bipartite networks (center) in which actors and groups appear as distinct kinds of nodes. Bipartite networks can always be projected onto one of two single-mode networks representing affiliations between the actors (bottom) or interlocks between the groups (top).

Hierarchical Networks

- P39, [Watts2003]

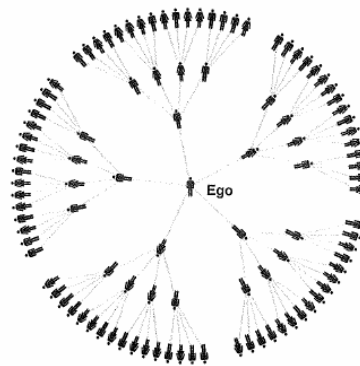


Figure 1.2. A pure branching network. Ego knows only 5 people, but within two degrees of separation, ego can reach 25; within three degrees, 105; and so on.

Formalizing the Small World Problem

[Watts 2003]

- Page 76 -82
- The alpha parameter

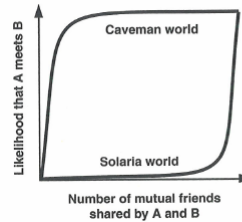


Figure 3.1. Two extreme kinds of interaction rules. In the top curve (caveman world), even a single mutual friend implies that A and B are highly likely to meet. In the bottom curve (Solaria world), all interactions are equally unlikely, regardless of how many friends A and B share.

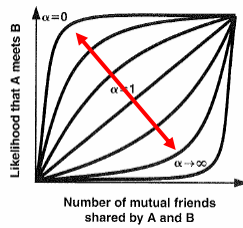


Figure 3.2. Between the two extremes, a whole range of interaction rules exists, specified by a particular value of the tuneable parameter alpha (α). When $\alpha = 0$, we have a caveman world; when α becomes infinite, we have Solaria.

Searchability

Two seemingly contradictory requirements for the Small World Phenomenon:

- Network should display a large clustering coefficient, so that a node's friends will know each other (as in Caveman world)
- It should be possible to connect two people chosen at random via chain of only a few intermediaries (as in Solaria world)

Formalizing the Small World Problem

[Watts 2003]

- Page 76 -82
- The alpha parameter
- Path length: computed only over nodes in the same connected component

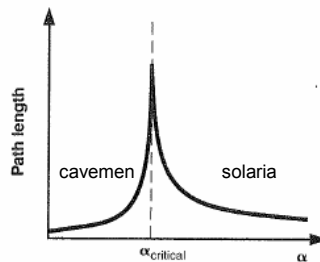


Figure 3.3. Path length as a function of alpha (α). At the critical alpha value, many small clusters join to connect the entire network, whose length then shrinks rapidly.

Formalizing the Small World Problem

- Page 76 -82
- Comparison between path length and clustering coefficient

Reminder - previous informal definition: SMP exists when every pair of nodes in a graph is connected by a path with an extremely small number of steps.
Does not take searchability into account. Random networks are hard to search with local knowledge

Small World Phenomenon exists when

$$L \geq L_{\text{random}} \text{ but } C \gg C_{\text{random}}$$

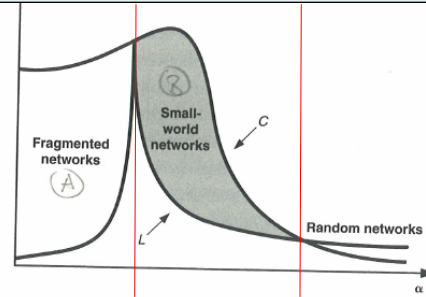


Figure 3.4. Comparison between path length (L) and clustering coefficient (C). The region between the curves, where L is small and C is large (shaded), represents the presence of small-world networks.

Examples for Small World Networks

[Watts and Strogatz 1998]

Table 1 Empirical examples of small-world networks

	L_{actual}	L_{random}	C_{actual}	C_{random}
Film actors	3.65	2.99	0.79	0.00027
Power grid	18.7	12.4	0.080	0.005
<i>C. elegans</i>	2.65	2.25	0.28	0.05

Characteristic path length L and clustering coefficient C for three real networks, compared to random graphs with the same number of vertices (n) and average number of edges per vertex (k). (Actors: $n = 225,226, k = 61$. Power grid: $n = 4,941, k = 2.67$. *C. elegans*: $n = 282, k = 14$.) The graphs are defined as follows. Two actors are joined by an edge if they have acted in a film together. We restrict attention to the giant connected component¹⁶ of this graph, which includes ~90% of all actors listed in the Internet Movie Database (available at <http://us.imdb.com>), as of April 1997. For the power grid, vertices represent generators, transformers and substations, and edges represent high-voltage transmission lines between them. For *C. elegans*, an edge joins two neurons if they are connected by either a synapse or a gap junction. We treat all edges as undirected and unweighted, and all vertices as identical, recognizing that these are crude approximations. **All three networks show the small-world phenomenon: $L \geq L_{\text{random}}$ but $C \gg C_{\text{random}}$.**

Any questions?

See you next week!